



TITLE:

Asymptotic distribution of the distance-based classifier under a strongly spiked eigenvalue model (Bayes Inference and Its Related Topics)

AUTHOR(S):

Ishii, Aki

CITATION:

Ishii, Aki. Asymptotic distribution of the distance-based classifier under a strongly spiked eigenvalue model (Bayes Inference and Its Related Topics). 数理解析研究所講究録 2017, 2047: 1-9

ISSUE DATE:

2017-10

URL:

<http://hdl.handle.net/2433/237024>

RIGHT:

Asymptotic distribution of the distance-based classifier under a strongly spiked eigenvalue model

東京理科大学・情報科学科 石井 晶 (Aki Ishii)
Department of Information Sciences
Tokyo University of Science

Abstract: We consider two-class classification for high-dimensional data. We consider the distance-based classifier given by Aoshima and Yata (2014). We provide an asymptotic distribution of the classifier under a strongly spiked eigenvalue model.

Key words and phrases: Asymptotic distribution, Distance-based classifier, HDLSS, Large p small n .

1. Introduction

Nowadays, you can see many types of high-dimensional data such as genetic microarrays, medical imaging, text recognition, finance, chemometrics, and so on. A common feature of high-dimensional data is that the data dimension is extremely high, however, the sample size is relatively low. We call such data “HDLSS” or “large p , small n ” data, where p is the data dimension and n is the sample size. In this paper, we consider two-class classification in HDLSS context. We aim to give an asymptotic distribution of the distance-based classifier under a strongly spiked eigenvalue model that was proposed by Aoshima and Yata (2017).

Suppose we have two classes π_i , $i = 1, 2$, and define independent $p \times n_i$ data matrices, $\mathbf{X}_i = [\mathbf{x}_{i1}, \dots, \mathbf{x}_{in_i}]$, $i = 1, 2$, from π_i , $i = 1, 2$, where \mathbf{x}_{ij} , $j = 1, \dots, n_i$, are independent and identically distributed (i.i.d.) as a p -dimensional distribution with a mean vector $\boldsymbol{\mu}_i$ and covariance matrix $\boldsymbol{\Sigma}_i (\geq \mathbf{O})$. We assume $n_i \geq 3$, $i = 1, 2$. The eigen-decomposition of $\boldsymbol{\Sigma}_i$ is given by

$$\boldsymbol{\Sigma}_i = \mathbf{H}_i \boldsymbol{\Lambda}_i \mathbf{H}_i^T = \sum_{s=1}^p \lambda_{s(i)} \mathbf{h}_{s(i)} \mathbf{h}_{s(i)}^T,$$

where $\boldsymbol{\Lambda}_i = \text{diag}(\lambda_{1(i)}, \dots, \lambda_{p(i)})$ having $\lambda_{1(i)} \geq \dots \geq \lambda_{p(i)} (\geq 0)$ and $\mathbf{H}_i = [\mathbf{h}_{1(i)}, \dots, \mathbf{h}_{p(i)}]$ is an orthogonal matrix of the corresponding eigenvectors. Let $\mathbf{X}_i - [\boldsymbol{\mu}_i, \dots, \boldsymbol{\mu}_i] = \mathbf{H}_i \boldsymbol{\Lambda}_i^{1/2} \mathbf{Z}_i$ for $i = 1, 2$. Then, \mathbf{Z}_i is a $p \times n_i$ sphered data matrix from a distribution with the zero mean and identity covariance matrix. Let $\mathbf{Z}_i = [\mathbf{z}_{1(i)}, \dots, \mathbf{z}_{p(i)}]^T$ and $\mathbf{z}_{j(i)} = (z_{j1(i)}, \dots, z_{jn_i(i)})^T$, $j = 1, \dots, p$, for $i = 1, 2$. Note that $E(z_{jk(i)} z_{j'k(i)}) = 0$ ($j \neq j'$) and $\text{Var}(\mathbf{z}_{j(i)}) = \mathbf{I}_{n_i}$, where \mathbf{I}_{n_i} denotes the n_i -dimensional identity matrix. Also, note that if \mathbf{X}_i is Gaussian, $z_{jk(i)}$ s are i.i.d. as the standard normal distribution, $N(0, 1)$. We assume that the fourth moments of each variable in \mathbf{Z}_i are uniformly bounded for $i = 1, 2$. Let $\mathbf{z}_{oj(i)} = \mathbf{z}_{j(i)} - (\bar{z}_{j(i)}, \dots, \bar{z}_{j(i)})^T$, $j = 1, \dots, p$; $i = 1, 2$, where $\bar{z}_{j(i)} = n_i^{-1} \sum_{k=1}^{n_i} z_{jk(i)}$.

We assume that $P(\lim_{p \rightarrow \infty} \|\mathbf{z}_{o1(i)}\| \neq 0) = 1$ for $i = 1, 2$, where $\|\cdot\|$ denotes the Euclidean norm.

Let \mathbf{x}_0 be an observation vector of an individual belonging to π_i ($i = 1, 2$). We assume \mathbf{x}_0 and \mathbf{x}_{ij} s are independent. We estimate $\boldsymbol{\mu}_i$ and $\boldsymbol{\Sigma}_i$ by $\bar{\mathbf{x}}_{in_i} = \sum_{j=1}^{n_i} \mathbf{x}_{ij}/n_i$ and $\mathbf{S}_{in_i} = \sum_{j=1}^{n_i} (\mathbf{x}_{ij} - \bar{\mathbf{x}}_{in_i})(\mathbf{x}_{ij} - \bar{\mathbf{x}}_{in_i})^T/(n_i - 1)$. A typical classification rule is that one classifies an individual into π_1 if

$$\begin{aligned} & (\mathbf{x}_0 - \bar{\mathbf{x}}_{1n_1})^T \mathbf{S}_{1n_1}^{-1} (\mathbf{x}_0 - \bar{\mathbf{x}}_{1n_1}) - \log \left\{ \frac{\det(\mathbf{S}_{2n_2})}{\det(\mathbf{S}_{1n_1})} \right\} \\ & < (\mathbf{x}_0 - \bar{\mathbf{x}}_{2n_2})^T \mathbf{S}_{2n_2}^{-1} (\mathbf{x}_0 - \bar{\mathbf{x}}_{2n_2}), \end{aligned} \quad (1.1)$$

and into π_2 otherwise. However, the inverse matrix of \mathbf{S}_{in_i} does not exist in the HDLSS context ($p > n_i$). When $\boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2$, Bickel and Levina (2004) considered the inverse matrix defined by only diagonal elements of the pooled sample covariance matrix. Yata and Aoshima (2012) considered using a ridge-type inverse covariance matrix derived by the *noise reduction (NR) methodology*. When $\boldsymbol{\Sigma}_1 \neq \boldsymbol{\Sigma}_2$, Dudoit et al. (2002) considered using the inverse matrix defined by only diagonal elements of \mathbf{S}_{in_i} . Aoshima and Yata (2011, 2015a) considered substituting $\{\text{tr}(\mathbf{S}_{in_i})/p\} \mathbf{I}_p$ for \mathbf{S}_{in_i} by using the difference of a geometric representation of HDLSS data from each π_i . Aoshima and Yata (2015b) considered quadratic classifiers in general and discussed asymptotic properties and optimality of the classifiers under high-dimensional settings. They showed that the misclassification rates tend to zero as the dimension goes to infinity. On the other hand, Hall et al. (2005, 2008), and Chan and Hall (2009) considered distance-based classifiers. Aoshima and Yata (2014) gave the misclassification rate adjusted classifier for multiclass, high-dimensional data whose misclassification rates are no more than specified thresholds under the following condition for eigenvalues:

$$\frac{\lambda_{1(i)}^2}{\text{tr}(\boldsymbol{\Sigma}_i^2)} \rightarrow 0 \text{ as } p \rightarrow \infty \text{ for } i = 1, 2. \quad (1.2)$$

Recently, Aoshima and Yata (2017) considered the “strongly spiked eigenvalue (SSE) model” as follows:

$$\liminf_{p \rightarrow \infty} \left\{ \frac{\lambda_{1(i)}^2}{\text{tr}(\boldsymbol{\Sigma}_i^2)} \right\} > 0 \text{ for } i = 1 \text{ or } 2. \quad (1.3)$$

On the other hand, Aoshima and Yata (2017) called (1.2) the “non-strongly spiked eigenvalue (NSSE) model”.

In this paper, we consider the distance-based classifier under one of the SSE models.

Remark 1.1. For a spiked model such as

$$\lambda_{s(i)} = a_{s(i)} p^{\alpha_{s(i)}} \quad (s = 1, \dots, t_i) \quad \text{and} \quad \lambda_{s(i)} = c_{s(i)} \quad (s = t_i + 1, \dots, p) \quad (1.4)$$

with positive and fixed constants, $a_{s(i)}$ s, $c_{s(i)}$ s and $\alpha_{s(i)}$ s, and a positive and fixed integer t_i , note that (1.2) holds when $\alpha_{1(i)} < 1/2$ for $i = 1, 2$. On the other hand, (1.3) holds for the spiked model in (1.4) with $\alpha_{1(i)} \geq 1/2$. See Yata and Aoshima (2012) for the details of the spiked model.

2. Distance-based classifier

Aoshima and Yata (2014) considered a classification rule given by using the identity matrix \mathbf{I}_p instead of \mathbf{S}_{in_i} in (1.1) as follows: One classifies an individual into π_1 if

$$\left(\mathbf{x}_0 - \frac{\bar{\mathbf{x}}_{1n_1} + \bar{\mathbf{x}}_{2n_2}}{2} \right)^T (\bar{\mathbf{x}}_{2n_2} - \bar{\mathbf{x}}_{1n_1}) - \frac{\text{tr}(\mathbf{S}_{1n_1})}{2n_1} + \frac{\text{tr}(\mathbf{S}_{2n_2})}{2n_2} < 0 \quad (2.1)$$

and into π_2 otherwise. Here, $-\text{tr}(\mathbf{S}_{1n_1})/(2n_1) + \text{tr}(\mathbf{S}_{2n_2})/(2n_2)$ is a bias-correction term. They showed the asymptotic normality of the classifier and provide a sample size determination so as to control misclassification rates being no more than a prespecified value. They further developed the classifier to multiclass classification.

Remark 2.1. Chan and Hall (2009) considered a scale adjusted distance-based classifier as follows: One classifies an individual into π_1 if

$$\begin{aligned} & \sum_{j=1}^{n_1} \frac{\|\mathbf{x}_0 - \mathbf{x}_{1j}\|^2}{n_1} - \sum_{j=1}^{n_2} \frac{\|\mathbf{x}_0 - \mathbf{x}_{2j}\|^2}{n_2} - \sum_{i=1}^{n_1} \sum_{j=1}^{n_1} \frac{\|\mathbf{x}_{1i} - \mathbf{x}_{1j}\|^2}{2n_1(n_1 - 1)} \\ & + \sum_{i=1}^{n_2} \sum_{j=1}^{n_2} \frac{\|\mathbf{x}_{2i} - \mathbf{x}_{2j}\|^2}{2n_2(n_2 - 1)} < 0 \end{aligned} \quad (2.2)$$

and into π_2 otherwise. We note that the classifier given by (2.1) is equivalent to the one given by (2.2), though the description of (2.1) is much simpler than (2.2).

We denote the error of misclassifying an individual from π_1 (into π_2) or π_2 (into π_1) by $e(1)$ or $e(2)$, respectively. Let $\Delta = \|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|^2$ and

$$W(\mathbf{x}_0) = \left(\mathbf{x}_0 - \frac{\bar{\mathbf{x}}_{1n_1} + \bar{\mathbf{x}}_{2n_2}}{2} \right)^T (\bar{\mathbf{x}}_{2n_2} - \bar{\mathbf{x}}_{1n_1}) - \frac{\text{tr}(\mathbf{S}_{1n_1})}{2n_1} + \frac{\text{tr}(\mathbf{S}_{2n_2})}{2n_2}.$$

Aoshima and Yata (2014) considered asymptotic properties of $W(\mathbf{x}_0)$ under the following assumptions:

(A-i) $\frac{(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}_i (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)}{\Delta^2} \rightarrow 0$ as $p \rightarrow \infty$ for $i = 1, 2$;

(A-ii) $\frac{\max_{j=1,2} \text{tr}(\boldsymbol{\Sigma}_j^2)}{n_i \Delta^2} \rightarrow 0$ as $p \rightarrow \infty$ either when n_i is fixed or $n_i \rightarrow \infty$ for $i = 1, 2$.

Then, they gave the asymptotic consistency:

Theorem 2.1 (Aoshima and Yata, 2014). Assume (A-i) and (A-ii). It holds that as $p \rightarrow \infty$

$$\frac{W(\mathbf{x}_0)}{\Delta} = \frac{(-1)^i}{2} + o_p(1) \quad \text{when } \mathbf{x}_0 \in \pi_i$$

for $i = 1, 2$. Then, the classification rule given by (2.1) has that as $p \rightarrow \infty$

$$e(1) \rightarrow 0 \quad \text{and} \quad e(2) \rightarrow 0.$$

Remark 2.2. Under the condition that $\max_{j=1,2}\{\text{tr}(\Sigma_j^2)\}/\Delta^2 \rightarrow 0$ as $p \rightarrow \infty$, one can claim Theorem 2.1 either when n_i is fixed or $n_i \rightarrow \infty$ for $i = 1, 2$.

Aoshima and Yata (2014) also showed the asymptotic normality of (2.1). They assume a general factor model as follows:

$$\mathbf{x}_{ij} = \Gamma_i \mathbf{w}_{ij} + \boldsymbol{\mu}_i$$

for $j = 1, \dots, n_i$; $i = 1, 2$, where Γ_i is a $p \times r_i$ matrix for some $r_i > 0$ such that $\Gamma_i \Gamma_i^T = \Sigma_i$, and \mathbf{w}_{ij} , $j = 1, \dots, n_i$, are i.i.d. random vectors having $E(\mathbf{w}_{ij}) = \mathbf{0}$ and $\text{Var}(\mathbf{w}_{ij}) = \mathbf{I}_{r_i}$. As for $\mathbf{w}_{ij} = (w_{i1j}, \dots, w_{ir_{ij}})^T$, $i = 1, 2$, we assume that

(A-iii) The fourth moments of each variable in \mathbf{w}_{ij} are uniformly bounded, $E(w_{iqj}^2 w_{isj}^2) = 1$ and $E(w_{iqj} w_{isj} w_{itj} w_{iuj}) = 0$ for all $q \neq s, t, u$.

If π_i is $N_p(\boldsymbol{\mu}_i, \Sigma_i)$, (A-iii) naturally follows. Also, Aoshima and Yata (2014) assume the following assumption for Σ_i , ($i = 1, 2$).

(A-iv) $\frac{\text{tr}(\Sigma_i \Sigma_l)}{\text{tr}(\Sigma_j^2)} \in (0, \infty)$ as $p \rightarrow \infty$ for $i, j, l = 1, 2$.

Here, $f(p) \in (0, \infty)$ as $p \rightarrow \infty$ denotes that $\liminf_{p \rightarrow \infty} f(p) > 0$ and $\limsup_{p \rightarrow \infty} f(p) < \infty$ for a function $f(\cdot)$. Let

$$\kappa_i = \frac{\text{tr}(\Sigma_i^2)}{n_i} + \frac{\text{tr}(\Sigma_1 \Sigma_2)}{n_j} + \sum_{i=1}^2 \frac{\text{tr}(\Sigma_i^2)}{2n_i(n_i - 1)}$$

for $i(\neq j) = 1, 2$. Let

$$n_{\min} = \min\{n_1, n_2\}.$$

We assume the following assumption:

(A-v) $\frac{(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \Sigma_i (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)}{\kappa_i} = o(1)$ as $p \rightarrow \infty$ and $n_{\min} \rightarrow \infty$ for $i = 1, 2$.

Then, they have the following result.

Theorem 2.2 (Aoshima and Yata, 2014). Assume (1.2). Assume also (A-iii) to (A-v). Then, we have that as $p \rightarrow \infty$ and $n_{\min} \rightarrow \infty$

$$\frac{W(\mathbf{x}_0) - (-1)^i \Delta/2}{\sqrt{\kappa_i}} \Rightarrow N(0, 1) \quad \text{when } \mathbf{x}_0 \in \pi_i \text{ for } i = 1, 2, \quad (2.3)$$

where “ \Rightarrow ” denotes the convergence in distribution and $N(0, 1)$ denotes a random variable distributed as the standard normal distribution.

Remark 2.3. From Theorem 2.2, for the classification rule by (2.1), it holds that as $p \rightarrow \infty$ and $n_{\min} \rightarrow \infty$

$$e(i) = \Phi\left(\frac{-\Delta}{2\sqrt{\kappa_i}}\right) + o(1) \quad \text{when } \mathbf{x}_0 \in \pi_i \text{ for } i = 1, 2 \quad (2.4)$$

under the assumptions in Theorem 2.2, where $\Phi(\cdot)$ denotes the cumulative distribution function of the standard normal distribution.

3. Asymptotic distribution under a SSE model

In this section, we provide the asymptotic distribution of (2.1) under a SSE model. Now, we consider the following assumptions for each π_i , $i = 1, 2$.

$$(B-i) \quad \frac{\text{tr}(\Sigma_i^2) - \lambda_{1(i)}^2}{\lambda_{1(i)}^2} = o(1) \quad \text{as } p \rightarrow \infty;$$

$$(B-ii) \quad \frac{\sum_{r,s \geq 2}^p \lambda_{r(i)} \lambda_{s(i)} E\{(z_{rk(i)}^2 - 1)(z_{sk(i)}^2 - 1)\}}{n_i \lambda_{1(i)}^2} = o(1) \quad \text{as } p \rightarrow \infty \text{ either when } n_i \text{ is fixed or } n_i \rightarrow \infty;$$

$$(B-iii) \quad z_{1k(i)}, k = 1, \dots, n_i \text{ i.i.d. as } N(0, 1).$$

Note that (B-i) is one of the SSE models. By using the NR method, $\lambda_{j(i)}$ s are estimated by

$$\tilde{\lambda}_{j(i)} = \hat{\lambda}_{j(i)} - \frac{\text{tr}(\mathbf{S}_{in_i}) - \sum_{s=1}^j \hat{\lambda}_{s(i)}}{n_i - 1 - j} \quad (j = 1, \dots, n_i - 2), \quad (3.1)$$

where $\hat{\lambda}_{j(i)}$ is the j th sample eigenvalue for $i = 1, 2$. Note that $\tilde{\lambda}_{j(i)} \geq 0$ w.p.1 for $j = 1, \dots, n_i - 2$. Yata and Aoshima (2012, 2013, 2016) showed that $\tilde{\lambda}_{j(i)}$ has several consistency properties when $p \rightarrow \infty$ and $n_i \rightarrow \infty$. On the other hand, when $p \rightarrow \infty$ while n_i s are fixed, Ishii et al. (2016) gave the following results.

Theorem 3.1 (Ishii et al. 2016). *Under (B-i) and (B-ii), it holds that as $p \rightarrow \infty$*

$$\frac{\tilde{\lambda}_{1(i)}}{\lambda_{1(i)}} = \begin{cases} \|z_{o1(i)}\|^2 / (n_i - 1) + o_p(1) & \text{when } n_i \text{ is fixed,} \\ 1 + o_p(1) & \text{when } n_i \rightarrow \infty \end{cases}$$

for $i = 1, 2$. Under (B-i) to (B-iii), it holds that as $p \rightarrow \infty$ when n_i is fixed

$$(n_i - 1) \frac{\tilde{\lambda}_{1(i)}}{\lambda_{1(i)}} \Rightarrow \chi_{n_i-1}^2 \quad \text{for } i = 1, 2.$$

Now, we consider the following assumptions.

$$(B-iv) \quad \frac{\lambda_{1(1)}}{\lambda_{1(2)}} = 1 + o(1) \text{ and } \mathbf{h}_{1(1)}^T \mathbf{h}_{1(2)} = 1 + o(1) \quad \text{as } p \rightarrow \infty.$$

$$(B-v) \quad \frac{(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \Sigma_i (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)}{\lambda_{1(i)}} = o(n_{\min}^{-1}) \quad \text{as } p \rightarrow \infty \text{ and } n_{\min} \rightarrow \infty \text{ for } i = 1, 2.$$

Note that (B-iv) means that the two class share their first eigenspace. One can check the validity of (B-iv) by using a test procedure given by Ishii et al. (2016).

Now, we consider the asymptotic distribution of (2.1) under the SSE model, (B-i). Let $z_{01(i)} = \mathbf{h}_{1(i)}^T(\mathbf{x}_0 - \boldsymbol{\mu}_i)/\lambda_{1(i)}^{1/2}$ when $\mathbf{x}_0 \in \pi_i$ for $i = 1, 2$. Then, we have the following result.

Lemma 3.1. *Assume (B-i), (B-iv) and (B-v). Then, we have that as $p \rightarrow \infty$ and $n_{\min} \rightarrow \infty$*

$$\frac{W(\mathbf{x}_0) - (-1)^i \Delta / 2}{\lambda_{1(i)}} = z_{01(i)}(\bar{z}_{1(2)} - \bar{z}_{1(1)}) + o_p(n_{\min}^{-1/2})$$

when $\mathbf{x}_0 \in \pi_i$ for $i = 1, 2$.

From Lemma 3.1 we have the following result.

Theorem 3.2. *Assume (B-i), (B-iii) to (B-v). Then, we have that as $p \rightarrow \infty$ and $n_{\min} \rightarrow \infty$*

$$u_n \frac{W(\mathbf{x}_0) - (-1)^i \Delta / 2}{\lambda_{1(i)}} \Rightarrow U_1 \times U_2$$

when $\mathbf{x}_0 \in \pi_i$ for $i = 1, 2$,

where $u_n = (n_1^{-1} + n_2^{-1})^{-1/2}$ and U_i s are mutually independent random variables distributed as $N(0, 1)$.

Remark 3.1. From Theorem 3.2, for the classification rule by (2.1), it holds that as $p \rightarrow \infty$ and $n_{\min} \rightarrow \infty$

$$e(i) = P \left\{ U_1 U_2 \leq -u_n \frac{\Delta}{2\lambda_{1(i)}} \right\} + o(1) \quad \text{when } \mathbf{x}_0 \in \pi_i \text{ for } i = 1, 2 \quad (3.2)$$

under the assumptions in Theorem 3.2. One can estimate Δ by

$$\hat{\Delta} = \|\bar{\mathbf{x}}_{1n_1} - \bar{\mathbf{x}}_{2n_2}\|^2 - \text{tr}(\mathbf{S}_{1n_1})/n_1 - \text{tr}(\mathbf{S}_{2n_2})/n_2.$$

The estimator was given by Aoshima and Yata (2011). Then, we can estimate (3.2) by $\hat{\Delta}$ and $\tilde{\lambda}_{1(i)}$.

Appendix

Proof of Lemma 3.1. We assume $\mathbf{x}_0 \in \pi_1$ without loss of generality. It holds that

$$\|\bar{\mathbf{x}}_{in_i} - \boldsymbol{\mu}_i\|^2 - \frac{\text{tr}(\mathbf{S}_{in_i})}{n_i} = \sum_{s=1}^p \sum_{k \neq k'} \lambda_{s(i)} \frac{z_{sk(i)} z_{sk'(i)}}{n_i(n_i - 1)}; \quad (A.1)$$

$$(\mathbf{x}_0 - \boldsymbol{\mu}_1)^T (\bar{\mathbf{x}}_{1n_1} - \boldsymbol{\mu}_1) = \lambda_{1(1)} z_{01(1)} \bar{z}_{1(1)} + \sum_{s=2}^p \lambda_{s(1)} z_{01(s)} \bar{z}_{s(1)}, \quad (A.2)$$

where $\mathbf{x}_0 - \boldsymbol{\mu}_1 = \mathbf{H}_1 \mathbf{\Lambda}_1^{1/2} (z_{01(1)}, \dots, z_{01(p)})^T$. By using Chebyshev's inequality, for any $\tau > 0$, under (B-i) and (B-iv), we have that as $p \rightarrow \infty$ and $n_{\min} \rightarrow \infty$

$$E \left\{ \left| \sum_{s=1}^p \sum_{k \neq k'} \frac{\lambda_{s(i)} z_{sk(i)} z_{sk'(i)}}{n_i(n_i - 1)} \right| \geq \tau n_{\min}^{-1/2} \lambda_{1(1)} \right\} \leq \frac{\sum_{s=1}^p \lambda_{s(i)}^2}{\tau^2 \lambda_{1(1)}^2 (n_i - 1)} \rightarrow 0;$$

$$E \left\{ \left| \sum_{s=2}^p \lambda_{s(1)} z_{0s(1)} \bar{z}_{s(1)} \right| \geq \tau n_{\min}^{-1/2} \lambda_{1(1)} \right\} \leq \frac{\sum_{s=2}^p \lambda_{s(1)}^2}{\tau^2 \lambda_{1(1)}^2} \rightarrow 0,$$

so that from (A.1) and (A.2)

$$\begin{aligned} \|\bar{\mathbf{x}}_{in_i} - \boldsymbol{\mu}_i\|^2 - \frac{\text{tr}(\mathbf{S}_{in_i})}{n_i} &= o_p(n_{\min}^{-1/2} \lambda_{1(1)}) \quad \text{for } i = 1, 2; \\ (\mathbf{x}_0 - \boldsymbol{\mu}_1)^T (\bar{\mathbf{x}}_{1n_1} - \boldsymbol{\mu}_1) &= \lambda_{1(1)} z_{01(1)} \bar{z}_{1(1)} + o_p(n_{\min}^{-1/2} \lambda_{1(1)}). \end{aligned} \quad (\text{A.3})$$

Let $\beta_{st} = (\lambda_{s(1)} \lambda_{t(2)})^{1/2} \mathbf{h}_{s(1)}^T \mathbf{h}_{t(2)}$ for all s, t . Then, we write that

$$\begin{aligned} (\mathbf{x}_0 - \boldsymbol{\mu}_1)^T (\bar{\mathbf{x}}_{2n_2} - \boldsymbol{\mu}_2) &= \sum_{s,t \geq 1}^p \beta_{st} z_{0s(1)} \bar{z}_{t(2)} \\ &= \beta_{11} z_{01(1)} \bar{z}_{1(2)} + \sum_{s=2}^p \beta_{s1} z_{0s(1)} \bar{z}_{1(2)} + \sum_{t=2}^p \beta_{1t} z_{01(1)} \bar{z}_{t(2)} \\ &\quad + \sum_{s,t \geq 2}^p \beta_{st} z_{0s(1)} \bar{z}_{t(2)}. \end{aligned} \quad (\text{A.4})$$

Let $\boldsymbol{\Sigma}_{i*} = \sum_{s=2}^p \lambda_{s(i)} \mathbf{h}_{s(i)} \mathbf{h}_{s(i)}^T$ for $i = 1, 2$. Here, we have that

$$\begin{aligned} E \left\{ \left(\sum_{s=2}^p \beta_{s1} z_{0s(1)} \bar{z}_{1(2)} \right)^2 \right\} &= \frac{\sum_{s=2}^p \beta_{s1}^2}{n_2} = \frac{\lambda_{1(2)} \mathbf{h}_{1(2)}^T \boldsymbol{\Sigma}_{1*} \mathbf{h}_{1(2)}}{n_2} \leq \frac{\lambda_{1(2)} \lambda_{2(1)}}{n_2}; \\ E \left\{ \left(\sum_{t=2}^p \beta_{1t} z_{01(1)} \bar{z}_{t(2)} \right)^2 \right\} &= \frac{\lambda_{1(1)} \mathbf{h}_{1(1)}^T \boldsymbol{\Sigma}_{2*} \mathbf{h}_{1(1)}}{n_2} \leq \frac{\lambda_{1(1)} \lambda_{2(2)}}{n_2}; \\ E \left\{ \left(\sum_{s,t \geq 2}^p \beta_{st} z_{0s(1)} \bar{z}_{t(2)} \right)^2 \right\} &= \frac{\text{tr}(\boldsymbol{\Sigma}_{1*} \boldsymbol{\Sigma}_{2*})}{n_2} \leq \frac{\sqrt{\text{tr}(\boldsymbol{\Sigma}_{1*}^2) \text{tr}(\boldsymbol{\Sigma}_{2*}^2)}}{n_2}. \end{aligned}$$

Then, by using Chebyshev's inequality, for any $\tau > 0$, under (B-i) and (B-iv), we have that as $p \rightarrow \infty$ and $n_{\min} \rightarrow \infty$

$$\begin{aligned} P \left(\left| \sum_{s=2}^p \beta_{s1} z_{0s(1)} \bar{z}_{1(2)} \right| > \tau n_{\min}^{-1/2} \lambda_{1(1)} \right) &\leq \frac{\lambda_{1(2)} \lambda_{2(1)}}{\tau^2 \lambda_{1(1)}^2} \rightarrow 0; \\ P \left(\left| \sum_{t=2}^p \beta_{1t} z_{01(1)} \bar{z}_{t(2)} \right| > \tau n_{\min}^{-1/2} \lambda_{1(1)} \right) &\leq \frac{\lambda_{1(1)} \lambda_{2(2)}}{\tau^2 \lambda_{1(1)}^2} \rightarrow 0; \\ P \left(\left| \sum_{s,t \geq 2}^p \beta_{st} z_{0s(1)} \bar{z}_{t(2)} \right| > \tau n_{\min}^{-1/2} \lambda_{1(1)} \right) &\leq \frac{\sqrt{\text{tr}(\boldsymbol{\Sigma}_{1*}^2) \text{tr}(\boldsymbol{\Sigma}_{2*}^2)}}{\tau^2 \lambda_{1(1)}^2} \rightarrow 0. \end{aligned}$$

Then, from (A.4) we have that

$$\begin{aligned} (\mathbf{x}_0 - \boldsymbol{\mu}_1)^T (\bar{\mathbf{x}}_{2n_2} - \boldsymbol{\mu}_2) &= \beta_{11} z_{01} \bar{z}_{1(2)} + o_p(\lambda_{1(1)} n_{\min}^{-1/2}) \\ &= \lambda_{1(1)} z_{01} \bar{z}_{1(2)} + o_p(\lambda_{1(1)} n_{\min}^{-1/2}). \end{aligned} \quad (\text{A.5})$$

Also, under (B-iv) and (B-v), we have that as $p \rightarrow \infty$ and $n_{\min} \rightarrow \infty$

$$\begin{aligned} P\left(\left|(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T (\mathbf{x}_0 - \boldsymbol{\mu}_1)\right| > \tau n_{\min}^{-1/2} \lambda_{1(1)}\right) &\leq \frac{n_{\min} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}_1 (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)}{\tau^2 \lambda_{1(1)}^2} \rightarrow 0; \\ P\left(\left|(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T (\bar{\mathbf{x}}_{2n_2} - \boldsymbol{\mu}_2)\right| > \tau n_{\min}^{-1/2} \lambda_{1(1)}\right) &\leq \frac{(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}_2 (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)}{\tau^2 \lambda_{1(1)}^2} \rightarrow 0, \end{aligned}$$

so that

$$\begin{aligned} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T (\mathbf{x}_0 - \boldsymbol{\mu}_1) &= o_p(n_{\min}^{-1/2} \lambda_{1(1)}); \\ (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T (\bar{\mathbf{x}}_{2n_2} - \boldsymbol{\mu}_2) &= o_p(n_{\min}^{-1/2} \lambda_{1(1)}) \end{aligned} \quad (\text{A.6})$$

Note that

$$\begin{aligned} W(\mathbf{x}_0) + \Delta/2 &= \frac{1}{2} \sum_{i=1}^2 (-1)^{i+1} \left\{ \|\bar{\mathbf{x}}_{in_i} - \boldsymbol{\mu}_i\|^2 - \frac{\text{tr}(\mathbf{S}_{in_i})}{n_i} \right\} \\ &\quad + \sum_{i=1}^2 (-1)^i (\mathbf{x}_0 - \boldsymbol{\mu}_1)^T (\bar{\mathbf{x}}_{in_i} - \boldsymbol{\mu}_i) \\ &\quad - (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T (\mathbf{x}_0 - \boldsymbol{\mu}_1) + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T (\bar{\mathbf{x}}_{2n_2} - \boldsymbol{\mu}_2). \end{aligned} \quad (\text{A.7})$$

Then, by combining (A.3), (A.5) and (A.6) with (A.7), under (B-i), (B-iv) and (B-v), it holds that as $p \rightarrow \infty$ and $n_{\min} \rightarrow \infty$

$$\frac{W(\mathbf{x}_0) + \Delta/2}{\lambda_{1(1)}} = z_{01(1)} (\bar{z}_{1(2)} - \bar{z}_{1(1)}) + o_p(n_{\min}^{-1/2}).$$

For the case when $\mathbf{x}_0 \in \pi_2$, we have the result similarly. Thus the proof is completed.

Proof of Theorem 3.2. By using Lemma 3.1, the result is obtained straightforwardly.

Acknowledgment

I would like to express my sincere gratitude to Professor Makoto Aoshima, for his enthusiastic guidance and helpful support to my research project. I would also like to thank Associate Professor, Kazuyoshi Yata, for his valuable suggestions.

References

- Aoshima, M., Yata, K. (2011). Two-stage procedures for high-dimensional data. *Sequential Analysis (Editor's special invited paper)*, 30, 356–399.
- Aoshima, M., Yata, K. (2014). A distance-based, misclassification rate adjusted classifier for multiclass, high-dimensional data. *Annals of the Institute of Statistical Mathematics*, 66, 983–1010.
- Aoshima, M., Yata, K. (2015a). Geometric classifier for multiclass, high-dimensional data. *Sequential Analysis*, 34, 279–294.
- Aoshima, M., Yata, K. (2015b). High-dimensional quadratic classifiers in non-sparse settings. *arXiv preprint*, arXiv:1503.04549.
- Aoshima, M., Yata, K. (2017). Two-sample tests for high-dimension, strongly spiked eigenvalue models. *Statistica Sinica*, in press (arXiv:1602.02491).
- Bickel, P.J., Levina, E. (2004). Some theory for Fisher's linear discriminant function, "naive Bayes", and some alternatives when there are many more variables than observations. *Bernoulli*, 10, 989–1010.
- Chan, Y.-B., Hall, P. (2009). Scale adjustments for classifiers in high-dimensional, low sample size settings. *Biometrika*, 96, 469–478.
- Dudoit, S., Fridlyand, J., Speed, T.P. (2002). Comparison of discrimination methods for the classification of tumors using gene expression data. *Journal of the American Statistical Association*, 97, 77–87.
- Hall, P., Marron, J.S., Neeman, A. (2005). Geometric representation of high dimension, low sample size data. *Journal of the Royal Statistical Society, Series B*, 67, 427–444.
- Hall, P., Pittelkow, Y., Ghosh, M. (2008). Theoretical measures of relative performance of classifiers for high dimensional data with small sample sizes. *Journal of the Royal Statistical Society, Series B*, 70, 159–173.
- Ishii, A., Yata, K., Aoshima, M. (2016). Asymptotic properties of the first principal component and equality tests of covariance matrices in high-dimension, low-sample-size context. *Journal of Statistical Planning and Inference*, 170, 186–199.
- Yata, K., Aoshima, M. (2012). Effective PCA for high-dimension, low-sample-size data with noise reduction via geometric representations. *Journal of Multivariate Analysis*, 105, 193–215.
- Yata, K., Aoshima, M. (2013). PCA consistency for the power spiked model in high-dimensional settings. *Journal of Multivariate Analysis*, 122, 334–354.
- Yata, K., Aoshima, M. (2016). Reconstruction of a high-dimensional low-rank matrix. *Electronic Journal of Statistics*, 10, 895–917.